

Mathmatial Foundations of Reinforcement Learning

Model-based RL with simulators



Yuting Wei

Statistics & Data Science, Wharton
University of Pennsylvania

Fall 2023

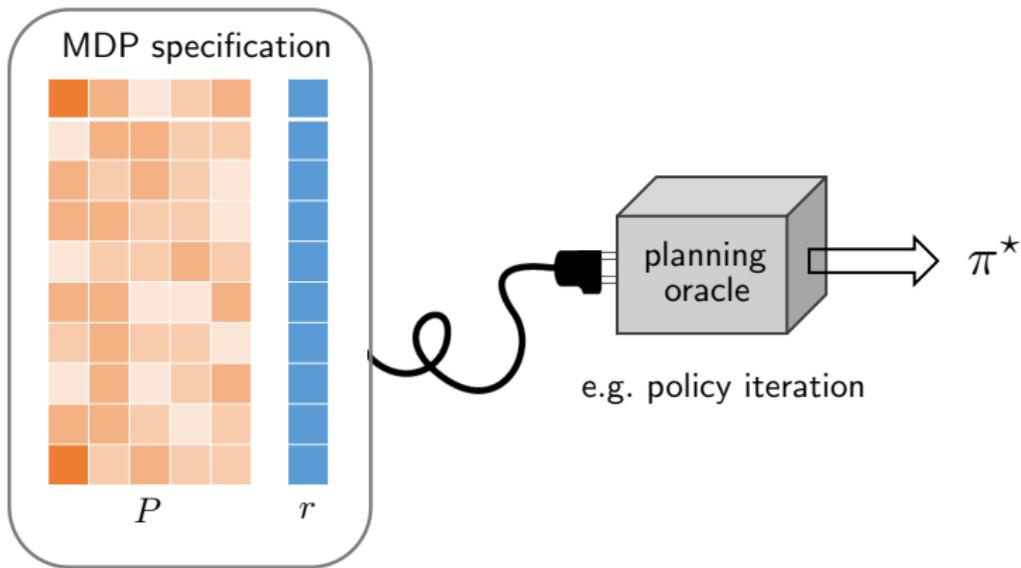
Outline

RL with a generative model

Model-based policy evaluation

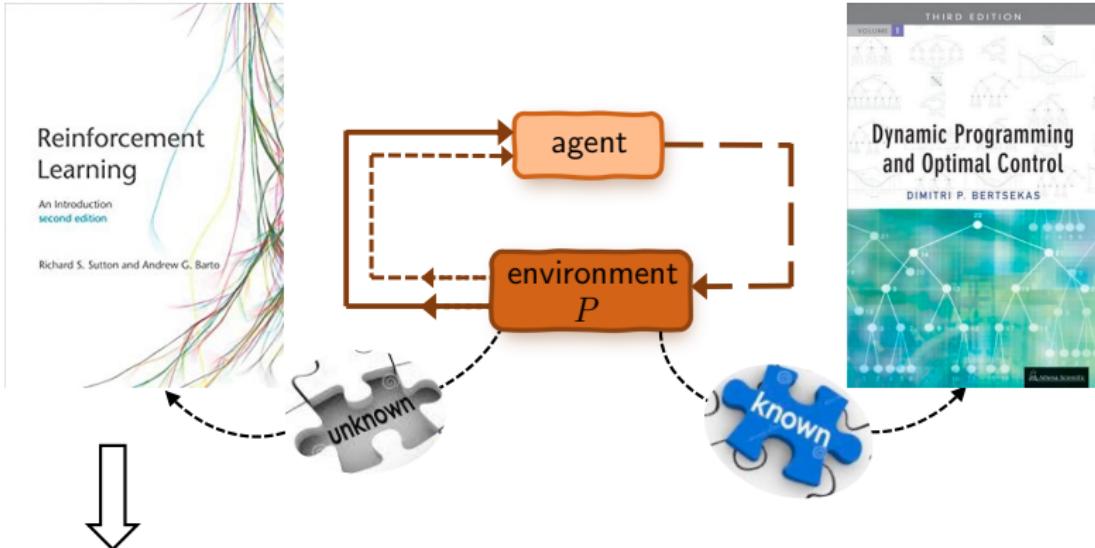
Model-based policy learning

Planning: when the MDP model is known



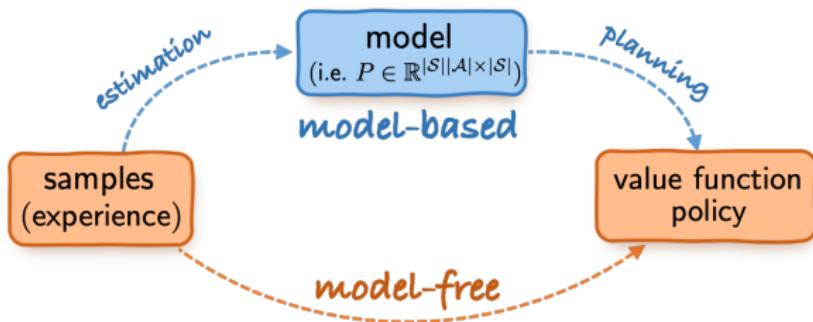
Planning: find the optimal policy π^* given MDP specification

Reinforcement learning (RL)



Learning: learn a desired policy from samples w/o model specification

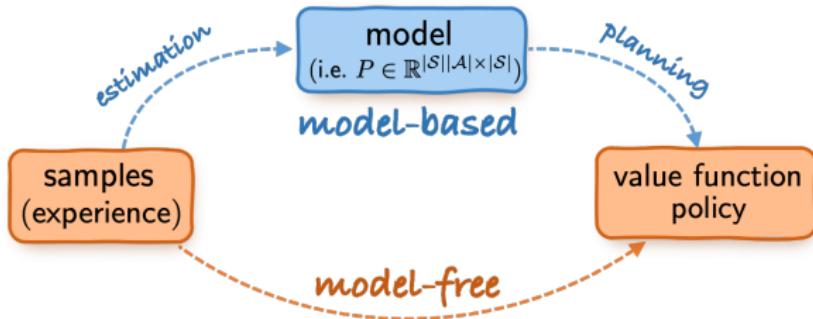
Two approaches to RL



Model-based approach (“plug-in”)

1. build an empirical estimate \hat{P} for P
2. planning based on empirical \hat{P}

Two approaches to RL



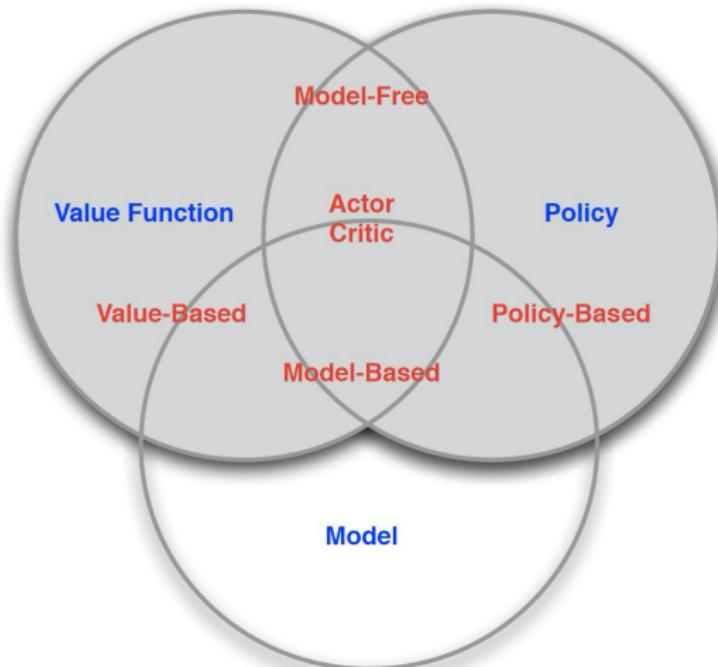
Model-based approach (“plug-in”)

1. build an empirical estimate \hat{P} for P
2. planning based on empirical \hat{P}

Model-free approach

— learning w/o constructing model explicitly

A taxonomy of RL approaches



—Credit: David Silver's slide

RL with a generative model

Motivation: study sample efficiency

Collecting data samples might be expensive or time-consuming



clinical trials



autonomous driving



online ads

Motivation: study sample efficiency

Collecting data samples might be expensive or time-consuming



clinical trials



autonomous driving



online ads

Understand and design of sample-efficient RL algorithms!

Data source in RL

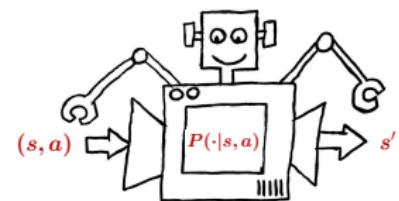
Exploration



offline RL



online RL



generative model

The capability of exploration increases from left to right.

Data source in RL

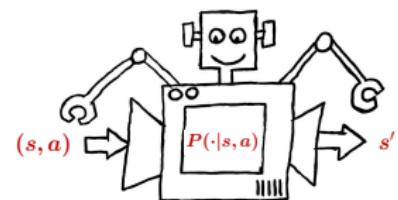
Exploration



offline RL



online RL



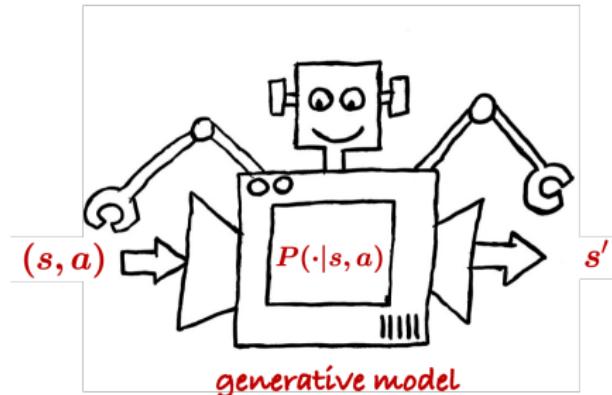
generative model

The capability of exploration increases from left to right.

This lecture: generative model / simulator

RL with a generative model / simulator

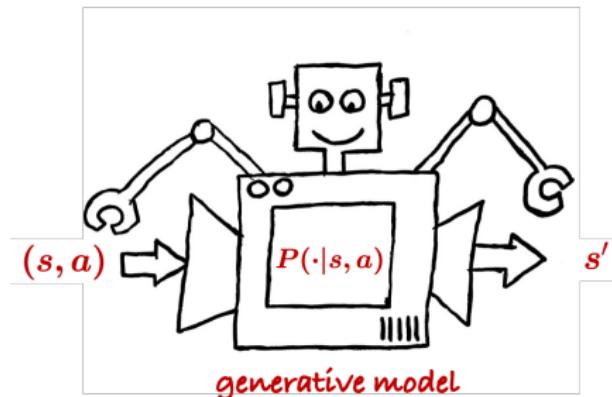
— [Kearns and Singh, 1999]



Protocol: for any state-action pair (s, a) , we can probe the simulator to output the next state s' .

RL with a generative model / simulator

— [Kearns and Singh, 1999]

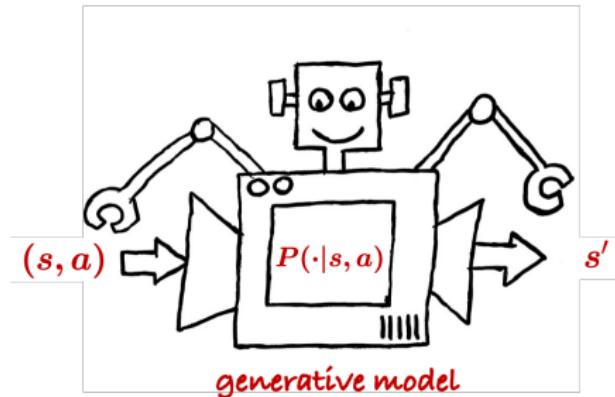


Protocol: for any state-action pair (s, a) , we can probe the simulator to output the next state s' .

We focus on the transition kernel and assume the reward is known or fixed, since the transition kernel captures the harder aspect of the problem.

RL with a generative model / simulator

— [Kearns and Singh, 1999]

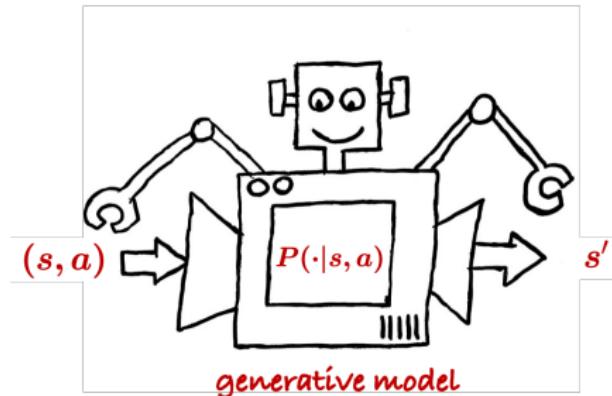


For each state-action pair (s, a) , collect N samples

$$\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$$

RL with a generative model / simulator

— [Kearns and Singh, 1999]

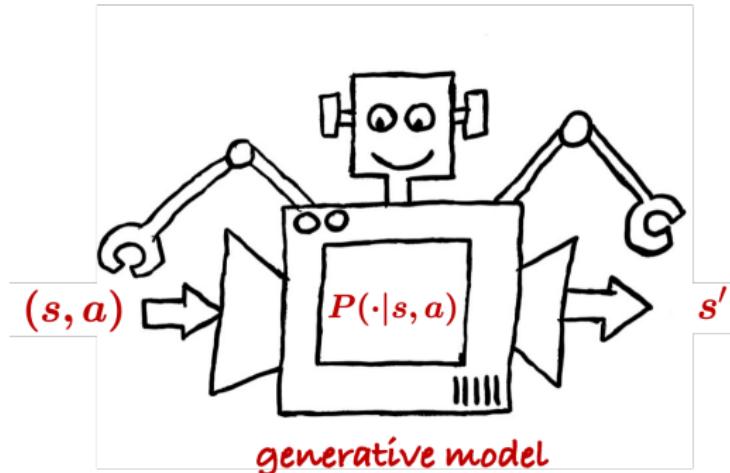


For each state-action pair (s, a) , collect N samples

$$\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$$

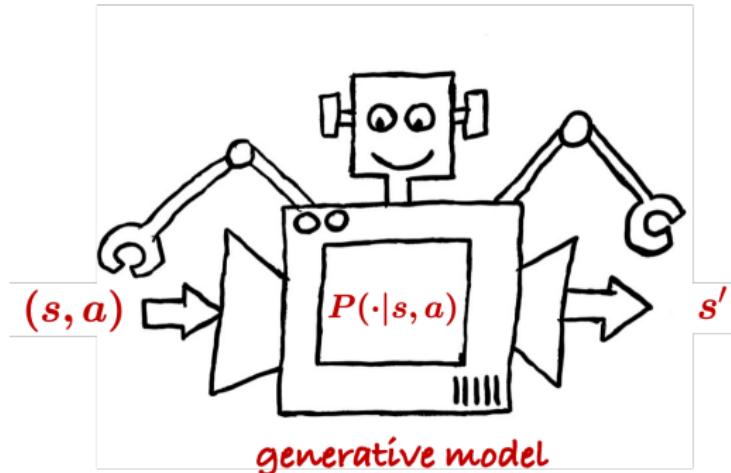
Question: How many samples are necessary and sufficient to solve the RL problem without worrying about exploration?

Model estimation under the generative model



For each (s, a) , collect N independent samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Model estimation under the generative model

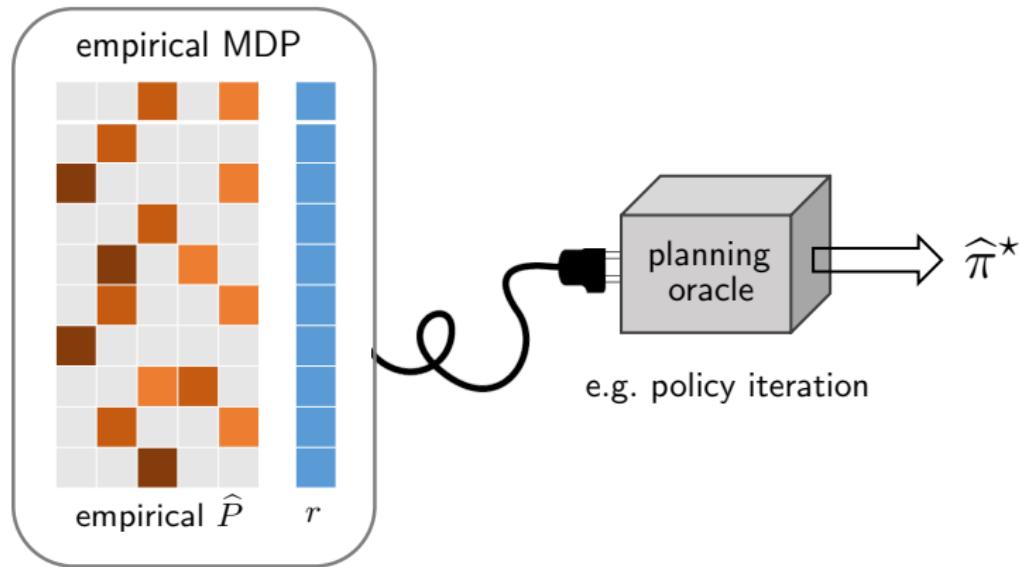


For each (s, a) , collect N independent samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Empirical estimates: estimate $\hat{P}(s'|s, a)$ by
$$\underbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s'_{(i)} = s'\}}_{\text{empirical frequency}}$$

Model-based (plug-in) estimator

— [Azar et al., 2013, Pananjady and Wainwright, 2020, Agarwal et al., 2020]



Run planning algorithms based on the *empirical* MDP

Questions

ℓ_∞ -sample complexity: how many samples are required for

- 1) evaluate an ε -accurate policy ?
- 2) learn an ε -optimal policy ?

Questions

ℓ_∞ -sample complexity: how many samples are required for

1) evaluate an ε -accurate policy ?

$$\forall s: |\widehat{V}^\pi(s) - V^\pi(s)| \leq \varepsilon$$

2) learn an ε -optimal policy ?

Questions

ℓ_∞ -sample complexity: how many samples are required for

1) evaluate an ε -accurate policy ?

$$\forall s: |\widehat{V}^\pi(s) - V^\pi(s)| \leq \varepsilon$$

2) learn an ε -optimal policy ?

$$\forall s: V^{\hat{\pi}}(s) \geq V^*(s) - \varepsilon$$

Model-based policy evaluation

Minimax lower bound

Theorem 1 (minimax lower bound; [Pananjady and Wainwright, 2020])

Fix a policy π . For all $\varepsilon \in [0, \frac{1}{1-\gamma})$, there exists some MDP such that the total number of samples need to be at least

$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2} \right)$$

to achieve $\|\hat{Q} - Q^\pi\|_\infty \leq \varepsilon$, where \hat{Q} is the output of any RL algorithm.

Minimax lower bound

Theorem 1 (minimax lower bound; [Pananjady and Wainwright, 2020])

Fix a policy π . For all $\varepsilon \in [0, \frac{1}{1-\gamma})$, there exists some MDP such that the total number of samples need to be at least

$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2} \right)$$

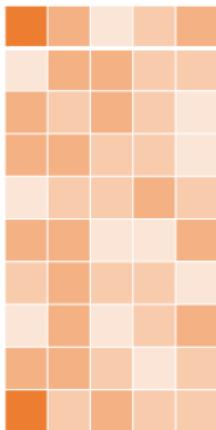
to achieve $\|\hat{Q} - Q^\pi\|_\infty \leq \varepsilon$, where \hat{Q} is the output of any RL algorithm.

- Consider the relative accuracy ε_{rel} by setting $\varepsilon := \frac{\varepsilon_{\text{rel}}}{1-\gamma}$, the lower bound can be equivalently expressed as

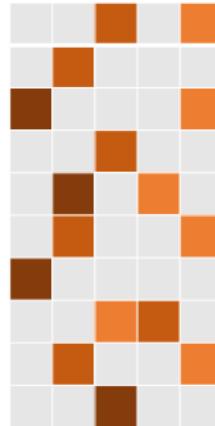
$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)\varepsilon_{\text{rel}}^2} \right).$$

- much smaller than the model dimension $|\mathcal{S}|^2|\mathcal{A}|$ — hint at the possibility of evaluating the policy without estimating the model reliably!

Challenges in the sample-starved regime



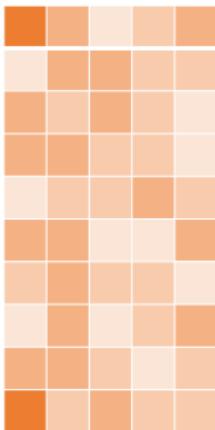
truth: $P \in \mathbb{R}^{|S||\mathcal{A}| \times |S|}$



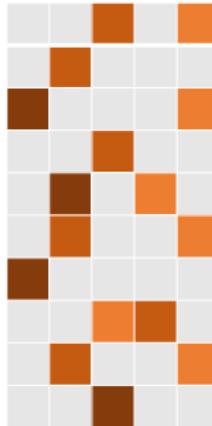
empirical estimate: \hat{P}

- Can't recover P faithfully if sample size $\ll |S|^2|\mathcal{A}|!$

Challenges in the sample-starved regime



truth: $P \in \mathbb{R}^{|S||\mathcal{A}| \times |S|}$



empirical estimate: \hat{P}

- Can't recover P faithfully if sample size $\ll |S|^2|\mathcal{A}|!$
- Can we trust our policy estimate when reliable model estimation is infeasible?

Recall: Bellman's consistency equation

- V^π / Q^π : value / action-value function under policy π

Bellman's consistency equation

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)]$$
$$Q^\pi(s, a) = \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\underbrace{V^\pi(s')}_{\text{next state's value}} \right]$$

The value/Q function can be decomposed into two parts:

- immediate reward $\mathbb{E}[r(s, a)]$
- discounted value of at the successor state
 $\gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} V(s')$



Richard Bellman

Policy evaluation for state-action function

Matrix-vector representation:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{\substack{s' \sim P(\cdot|s,a), a' \sim \pi(\cdot|s')}} [Q^\pi(s', a')]$$

⇓

$$Q^\pi = r + \gamma P^\pi Q^\pi$$

⇓

$$Q^\pi = (I - \gamma P^\pi)^{-1} r$$

- Here, P^π is the state-action transition matrix induced by π , namely,

$$P^\pi(s', a'|s, a) = P(s'|s, a)\pi(a'|s').$$

Sample complexity for plug-in policy evaluation

Model-based plug-in estimate:

$$\widehat{Q}^\pi = (I - \gamma \widehat{P}^\pi)^{-1} r$$

Theorem 2 ([Pananjady and Wainwright, 2020])

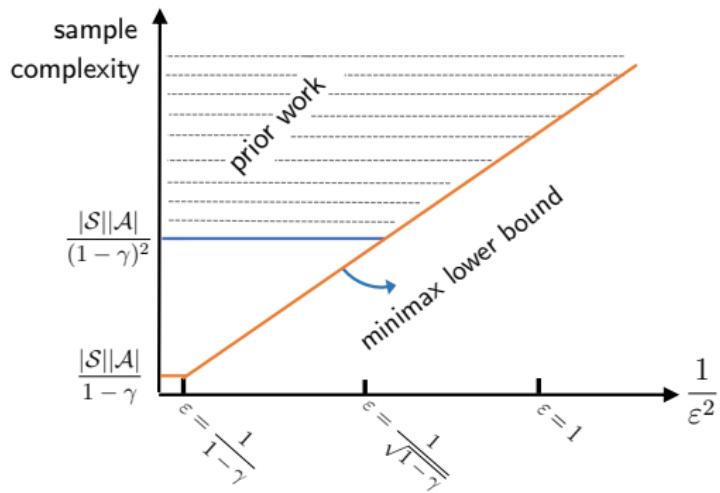
Fix any policy π . For $0 < \varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$, the plug-in estimator \widehat{Q}^π obeys

$$\|\widehat{Q}^\pi - Q^\pi\|_\infty \leq \varepsilon$$

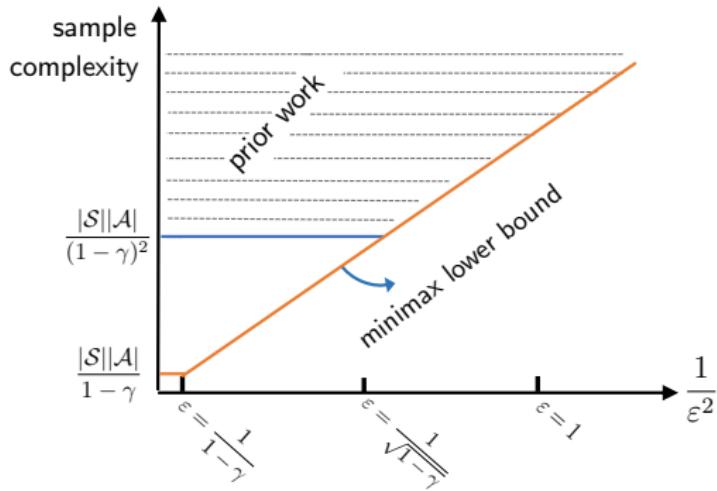
with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2}\right).$$

A sample size barrier



A sample size barrier



- A sample size barrier $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2}$ appeared in prior works [Azar et al., 2013, Pananjady and Wainwright, 2020]

Refined analysis

Model-based plug-in estimate:

$$\hat{Q}^\pi = (I - \gamma \hat{P}^\pi)^{-1} r$$

Theorem 3 ([Li et al., 2020])

Fix any policy π . For $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the plug-in estimator \hat{Q}^π obeys

$$\|\hat{Q}^\pi - Q^\pi\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2}\right)$$

Refined analysis

Model-based plug-in estimate:

$$\hat{Q}^\pi = (I - \gamma \hat{P}^\pi)^{-1} r$$

Theorem 3 ([Li et al., 2020])

Fix any policy π . For $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the plug-in estimator \hat{Q}^π obeys

$$\|\hat{Q}^\pi - Q^\pi\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2}\right)$$

- Minimax optimal for all ε [Azar et al., 2013, Pananjady and Wainwright, 2020].

Analysis: crude idea

- We'll demonstrate a crude version based on Hoeffding's inequality.

$$Q^\pi = (I - \gamma P^\pi)^{-1} r, \quad \hat{Q}^\pi = (I - \gamma \hat{P}^\pi)^{-1} r$$

Useful expansion:

$$\begin{aligned}\hat{Q}^\pi - Q^\pi &= (I - \gamma \hat{P}^\pi)^{-1} r - (I - \gamma P^\pi)^{-1} r \\&= (I - \gamma \hat{P}^\pi)^{-1} \left((I - \gamma P^\pi) - (I - \gamma \hat{P}^\pi) \right) Q^\pi \\&= \gamma (I - \gamma \hat{P}^\pi)^{-1} (\hat{P}^\pi - P^\pi) Q^\pi \\&= \gamma (I - \gamma \hat{P}^\pi)^{-1} (\hat{P} - P) V^\pi.\end{aligned}$$

Analysis: Hoeffding's inequality

By Hoeffding's inequality and union bound, with probability at least $1 - \delta$,

$$\left\| (\hat{P} - P)V^\pi \right\|_\infty \leq \sqrt{\frac{2 \log(2|\mathcal{S}||\mathcal{A}|/\delta)}{N(1-\gamma)^2}}.$$

Then, using $(I - \gamma \hat{P}^\pi)^{-1} = \sum_{i=0}^{\infty} \gamma^i (\hat{P}^\pi)^i$,

$$\begin{aligned} \left\| \gamma(I - \gamma \hat{P}^\pi)^{-1}(\hat{P} - P)V^\pi \right\|_\infty &\leq \gamma \sum_{i=0}^{\infty} \gamma^i \left\| (\hat{P}^\pi)^i (\hat{P} - P)V^\pi \right\|_\infty \\ &\leq \gamma \sum_{i=0}^{\infty} \gamma^i \left\| (\hat{P} - P)V^\pi \right\|_\infty \\ &\leq \frac{\gamma}{1-\gamma} \left\| (\hat{P} - P)V^\pi \right\|_\infty \leq \gamma \sqrt{\frac{2 \log(2|\mathcal{S}||\mathcal{A}|/\delta)}{N(1-\gamma)^4}}. \end{aligned}$$

Analysis: variance control + a peeling argument

- Better concentration with variance control: Bernstein's inequality
- Going beyond the 1st-order error expansion

$$\widehat{Q}^\pi - Q^\pi = \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)\widehat{Q}^\pi$$

Instead: higher-order expansion \longrightarrow tighter control

$$\widehat{Q}^\pi - Q^\pi = \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)Q^\pi +$$

Analysis: variance control + a peeling argument

- Better concentration with variance control: Bernstein's inequality
- Going beyond the 1st-order error expansion

$$\widehat{Q}^\pi - Q^\pi = \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)\widehat{Q}^\pi$$

Instead: higher-order expansion \longrightarrow tighter control

$$\begin{aligned}\widehat{Q}^\pi - Q^\pi &= \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)Q^\pi + \\ &\quad + \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)(\widehat{Q}^\pi - Q^\pi)\end{aligned}$$

Analysis: variance control + a peeling argument

- Better concentration with variance control: Bernstein's inequality
- Going beyond the 1st-order error expansion

$$\widehat{Q}^\pi - Q^\pi = \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)\widehat{Q}^\pi$$

Instead: higher-order expansion \longrightarrow tighter control

$$\begin{aligned}\widehat{Q}^\pi - Q^\pi &= \gamma(I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi)Q^\pi + \\ &\quad + \gamma^2 \left((I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi) \right)^2 Q^\pi \\ &\quad + \gamma^3 \left((I - \gamma P^\pi)^{-1}(\widehat{P}^\pi - P^\pi) \right)^3 Q^\pi \\ &\quad + \dots\end{aligned}$$

Model-based policy learning

Minimax lower bound

Theorem 4 (minimax lower bound; [Azar et al., 2013])

For all $\varepsilon \in [0, \frac{1}{1-\gamma})$, there exists some MDP such that the total number of samples need to be at least

$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2} \right)$$

to achieve $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$, where \hat{Q} is the output of any RL algorithm.

Minimax lower bound

Theorem 4 (minimax lower bound; [Azar et al., 2013])

For all $\varepsilon \in [0, \frac{1}{1-\gamma})$, there exists some MDP such that the total number of samples need to be at least

$$\tilde{\Omega} \left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \varepsilon^2} \right)$$

to achieve $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$, where \hat{Q} is the output of any RL algorithm.

- holds for both value-based and policy-based algorithms.
- much smaller than the model dimension $|\mathcal{S}|^2 |\mathcal{A}|$ — hint at the possibility of solving RL without estimating the model reliably!

Sample complexity for learning Q^*

Theorem 5 ([Azar et al., 2013])

For any $0 < \varepsilon \leq 1$, the optimal Q -function \widehat{Q} of the empirical MDP achieves

$$\|\widehat{Q} - Q^*\|_\infty \leq \varepsilon$$

with sample complexity at most $\widetilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$.

- matches with the minimax lower bound whenever $\varepsilon \in (0, 1]$.

Sample complexity for learning Q^*

Theorem 5 ([Azar et al., 2013])

For any $0 < \varepsilon \leq 1$, the optimal Q -function \widehat{Q} of the empirical MDP achieves

$$\|\widehat{Q} - Q^*\|_\infty \leq \varepsilon$$

with sample complexity at most $\widetilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$.

- matches with the minimax lower bound whenever $\varepsilon \in (0, 1]$.
- **Question:** Does it imply a near minimax-optimal policy $\widehat{\pi}$?

From Q-function to policy

Lemma 6 ([Singh and Yee, 1994])

Let the greedy policy w.r.t. \widehat{Q} be $\widehat{\pi}$, then

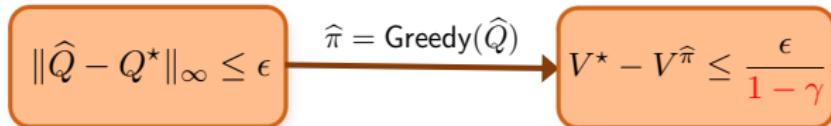
$$V^* - V^{\widehat{\pi}} \leq \frac{2}{1-\gamma} \|Q^* - \widehat{Q}\|_\infty.$$

From Q-function to policy

Lemma 6 ([Singh and Yee, 1994])

Let the greedy policy w.r.t. \hat{Q} be $\hat{\pi}$, then

$$V^* - V^{\hat{\pi}} \leq \frac{2}{1-\gamma} \|Q^* - \hat{Q}\|_\infty.$$

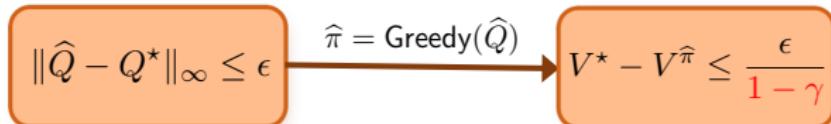


From Q-function to policy

Lemma 6 ([Singh and Yee, 1994])

Let the greedy policy w.r.t. \widehat{Q} be $\widehat{\pi}$, then

$$V^* - V^{\widehat{\pi}} \leq \frac{2}{1-\gamma} \|Q^* - \widehat{Q}\|_\infty.$$



This **error amplification** has consequences in sample complexities.

- To reach ϵ -optimality, the greedy policy of a minimax-optimal Q-function estimator needs

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^5 \varepsilon^2}\right)$$

samples invoking the above naive argument. Need refined arguments!

Theory of model-based policy learning

Theorem 7 ([Agarwal et al., 2020])

For any $0 < \varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$, the optimal policy $\hat{\pi}^*$ of the empirical MDP achieves

$$\|V^{\hat{\pi}^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

Theory of model-based policy learning

Theorem 7 ([Agarwal et al., 2020])

For any $0 < \varepsilon \leq \frac{1}{\sqrt{1-\gamma}}$, the optimal policy $\hat{\pi}^*$ of the empirical MDP achieves

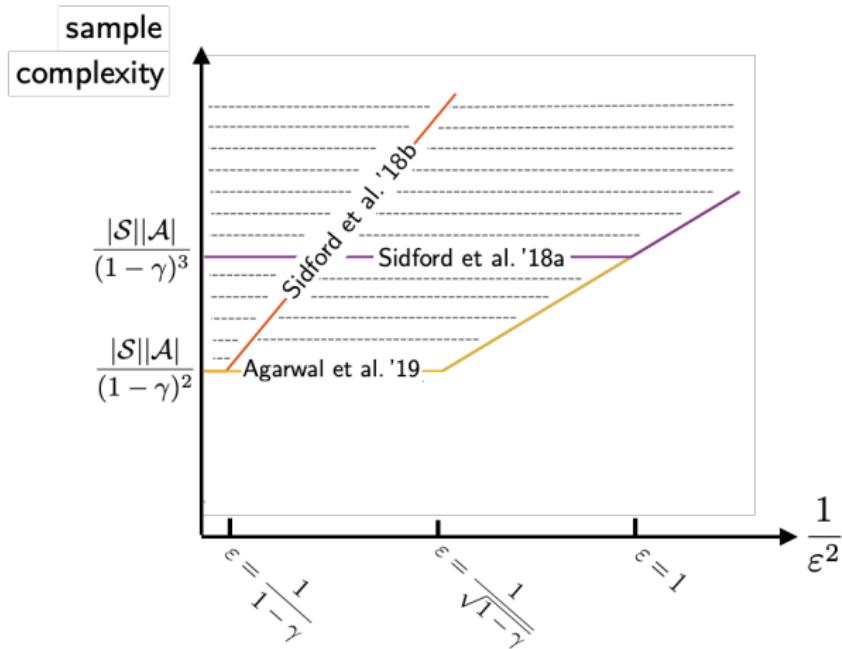
$$\|V^{\hat{\pi}^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

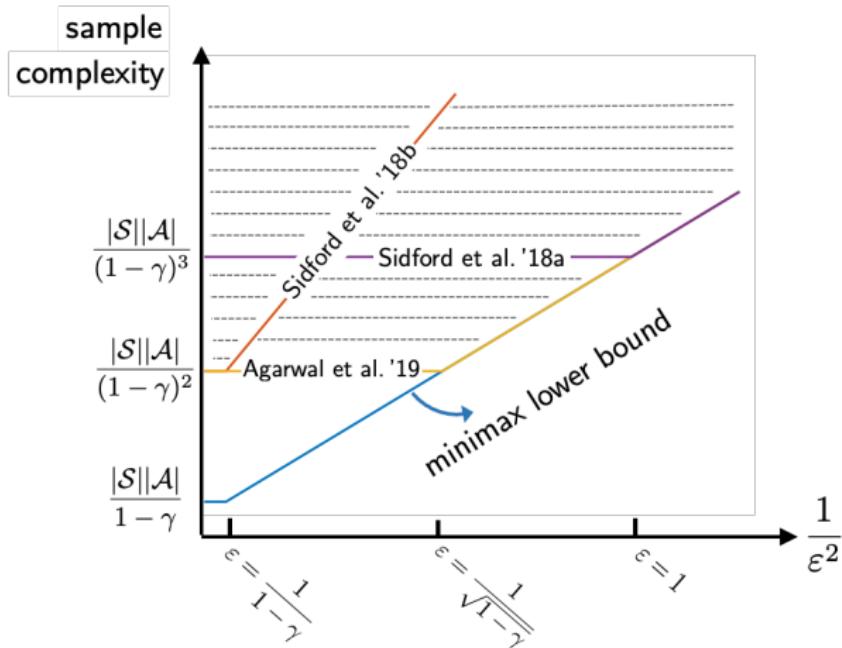
$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- Matches with the lower bound $\tilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ [Azar et al., 2013] when $\varepsilon \in (0, \frac{1}{\sqrt{1-\gamma}}]$.

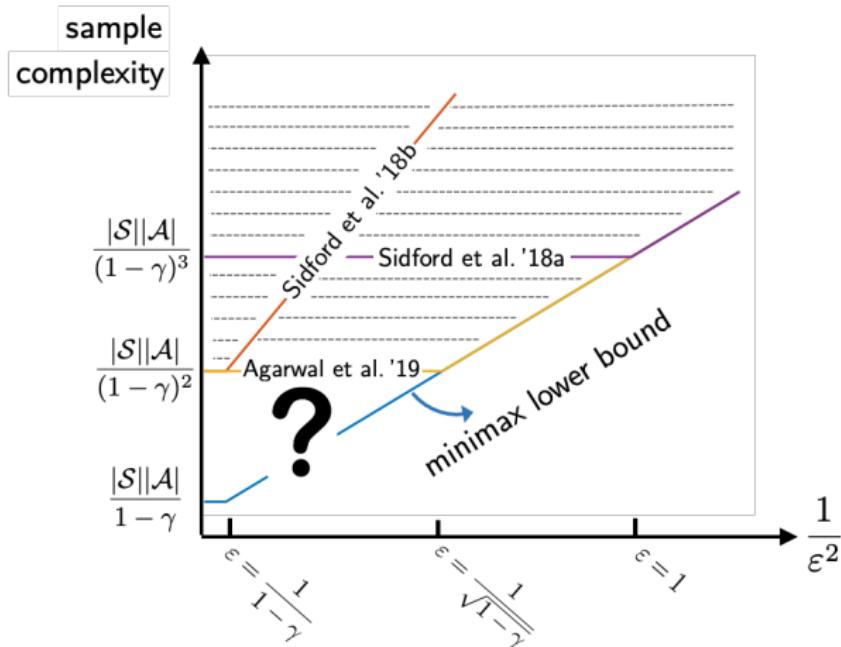
A sample complexity barrier



A sample complexity barrier



A sample complexity barrier

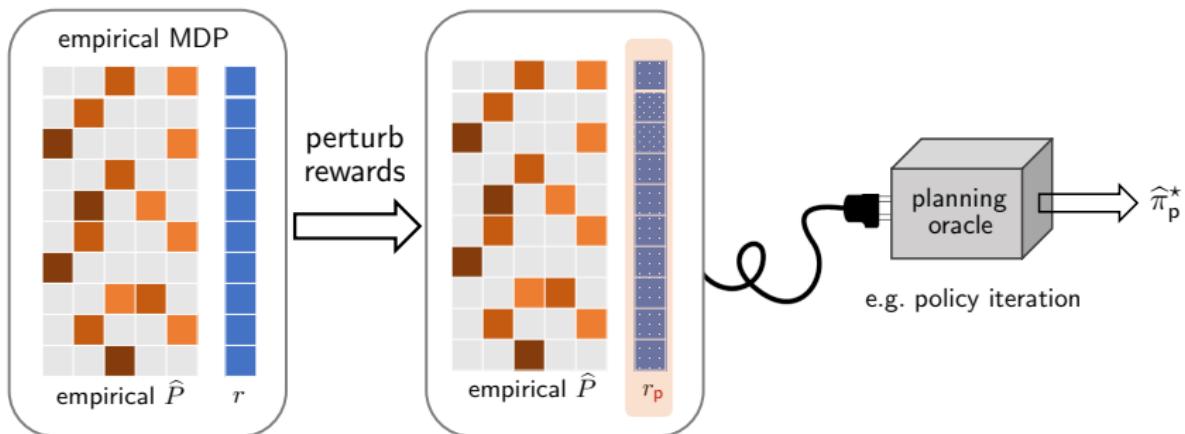


All prior theory requires sample size $> \underbrace{\frac{|S||\mathcal{A}|}{(1-\gamma)^2}}_{\text{sample size barrier}}$

Is it possible to close the gap?

Model-based plug-in estimator + perturbation

— [Li et al., 2020]



Planning based on the *empirical* MDP with *slightly perturbed rewards*

Refined theory of model-based policy learning

Theorem 8 ([Li et al., 2020])

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of the perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

Refined theory of model-based policy learning

Theorem 8 ([Li et al., 2020])

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of the perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- $\widehat{\pi}_p^*$: obtained by empirical VI or PI within $\tilde{O}\left(\frac{1}{1-\gamma}\right)$ iterations

Refined theory of model-based policy learning

Theorem 8 ([Li et al., 2020])

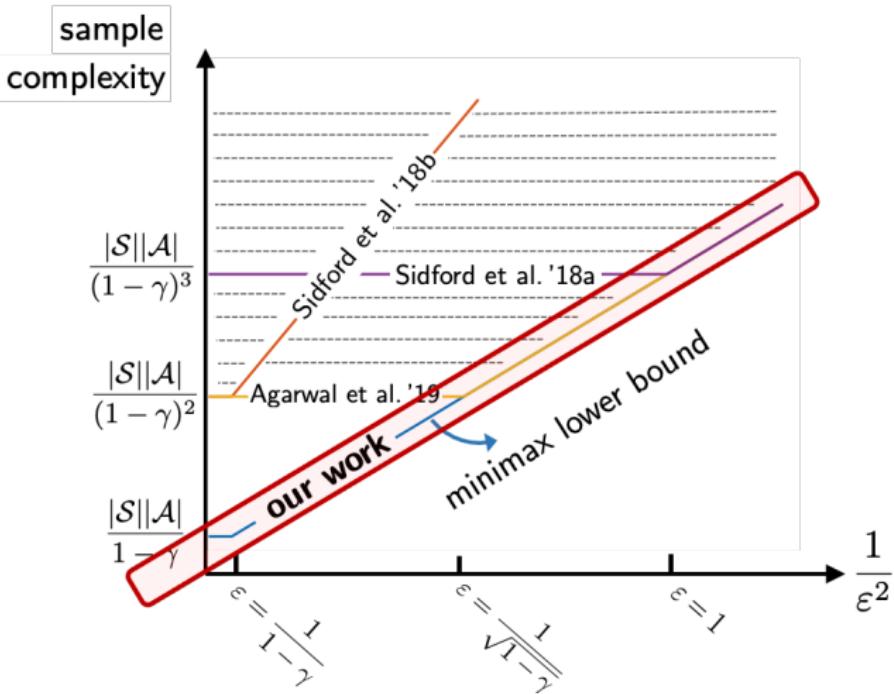
For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of the perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- $\widehat{\pi}_p^*$: obtained by empirical VI or PI within $\tilde{O}\left(\frac{1}{1-\gamma}\right)$ iterations
- **Minimax lower bound:** $\widetilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ [Azar et al., 2013]



Model-based RL is nearly minimax optimal!

Notation and Bellman equation

- V^π : true value function under policy π
 - Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1} r$

Notation and Bellman equation

- V^π : true value function under policy π
 - Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1} r$
- \hat{V}^π : estimate of value function under policy π
 - Bellman equation: $\hat{V}^\pi = (I - \gamma \hat{P}_\pi)^{-1} r$

Notation and Bellman equation

- V^π : true value function under policy π
 - Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1} r$
- \hat{V}^π : estimate of value function under policy π
 - Bellman equation: $\hat{V}^\pi = (I - \gamma \hat{P}_\pi)^{-1} r$
- π^* : optimal policy w.r.t. true value function
- $\hat{\pi}^*$: optimal policy w.r.t. empirical value function

Notation and Bellman equation

- V^π : true value function under policy π
 - Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1}r$
- \hat{V}^π : estimate of value function under policy π
 - Bellman equation: $\hat{V}^\pi = (I - \gamma \hat{P}_\pi)^{-1}r$
- π^* : optimal policy w.r.t. true value function
- $\hat{\pi}^*$: optimal policy w.r.t. empirical value function
- $V^* := V^{\pi^*}$: optimal values under true models
- $\hat{V}^* := \hat{V}^{\hat{\pi}^*}$: optimal values under empirical models

Proof ideas

Elementary decomposition:

$$V^* - V^{\widehat{\pi}^*} = (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{V}^{\widehat{\pi}^*}) + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*})$$

Proof ideas

Elementary decomposition:

$$\begin{aligned} V^* - \widehat{V}^{\pi^*} &= (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{\widehat{V}}^{\pi^*}) + (\widehat{\widehat{V}}^{\pi^*} - \widehat{V}^{\pi^*}) \\ &\leq (V^{\pi^*} - \widehat{V}^{\pi^*}) + \mathbf{0} + (\widehat{\widehat{V}}^{\pi^*} - \widehat{V}^{\pi^*}) \end{aligned}$$

Proof ideas

Elementary decomposition:

$$\begin{aligned} V^* - \widehat{V}^{\pi^*} &= (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{\widehat{V}}^{\pi^*}) + (\widehat{\widehat{V}}^{\pi^*} - \widehat{V}^{\pi^*}) \\ &\leq (V^{\pi^*} - \widehat{V}^{\pi^*}) + \textcolor{red}{0} + (\widehat{\widehat{V}}^{\pi^*} - \widehat{V}^{\pi^*}) \end{aligned}$$

- **Step 1:** control $V^\pi - \widehat{V}^\pi$ for a fixed π
(Bernstein inequality + high-order decomposition)

Proof ideas

Elementary decomposition:

$$\begin{aligned} V^* - \widehat{V}^{\pi^*} &= (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{V}^{\widehat{\pi}^*}) + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}) \\ &\leq (V^{\pi^*} - \widehat{V}^{\pi^*}) + \textcolor{red}{0} + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}) \end{aligned}$$

- **Step 1:** control $V^\pi - \widehat{V}^\pi$ for a fixed π
(Bernstein inequality + high-order decomposition)
- **Step 2:** extend it to control $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$ ($\widehat{\pi}^*$ depends on samples)
(decouple statistical dependency via leave-one-out analysis and reward perturbation)

References I

-  Agarwal, A., Kakade, S., and Yang, L. F. (2020).
Model-based reinforcement learning with a generative model is minimax optimal.
Conference on Learning Theory, pages 67–83.
-  Azar, M. G., Munos, R., and Kappen, H. J. (2013).
Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model.
Machine learning, 91(3):325–349.
-  Kearns, M. J. and Singh, S. P. (1999).
Finite-sample convergence rates for Q-learning and indirect algorithms.
In *Advances in neural information processing systems*, pages 996–1002.
-  Li, G., Wei, Y., Chi, Y., and Chen, Y. (2020).
Breaking the sample size barrier in model-based reinforcement learning with a generative model.
arXiv preprint arXiv:2005.12900.

References II



Pananjady, A. and Wainwright, M. J. (2020).

Instance-dependent ℓ_∞ -bounds for policy evaluation in tabular reinforcement learning.

IEEE Transactions on Information Theory, 67(1):566–585.



Singh, S. P. and Yee, R. C. (1994).

An upper bound on the loss from approximate optimal-value functions.

Machine Learning, 16:227–233.